# K-StyleLoRA: Information-Guided Image Generation via Selective Feature Learning

Anonymous ICCV submission

Paper ID *****

## Abstract

*Despite remarkable advances in image generation, existing diffusion models struggle to capture diverse cultural aesthetics. While Low-Rank Adaptation (LoRA) enables efficient fine-tuning, conventional approaches lack semantic awareness and apply uniform adaptations across all features, leading to suboptimal cultural representation. To address these limitations, we introduce **K-StyleLoRA**, a novel framework that leverages CLIP's cross-modal understanding for culturally-aware image generation. Our approach consists of two key innovations. First, CLIP-Guided Information Gating dynamically modulates LoRA adaptations based on cultural relevance scores, enabling selective enhancement of culturally-relevant features while suppressing irrelevant ones. Second, Cultural Semantic Loss provides additional semantic guidance through CLIP-based similarity optimization with Korean cultural concepts. Extensive experiments on Korean traditional art demonstrate superior cultural fidelity while maintaining generation quality and diversity. Most notably, K-StyleLoRA demonstrates exceptional cultural transfer capability on generic prompts requiring implicit cultural understanding, achieving a Cultural Similarity Score of 0.274, representing a 9.6% improvement over the vanilla SDXL baseline (0.250). Our framework establishes semantic-aware adaptation as a powerful paradigm for cultural representation, offering a scalable approach that can be extended to diverse cultural contexts and generation tasks beyond Korean aesthetics[1].*

## 1. Introduction

Image generation has undergone remarkable advancement with the emergence of large-scale diffusion models [1–3]. These models demonstrate exceptional capability in producing high-quality, diverse images from various input modalities. However, a critical limitation persists: existing models exhibit pronounced cultural bias, predominantly reflecting the cultural contexts most prevalent in their training data while struggling to authentically represent diverse cultural traditions [4, 5]. This bias stems from training data imbalances and the inherent difficulty of capturing nuanced cultural semantics across different generation tasks.

This cultural limitation becomes particularly evident when examining culturally-specific concepts. Consider the prompt *"traditional clothes"* across different cultural contexts: Korean hanbok with vibrant colors and flowing lines, Indian lehenga with intricate embroidery, Japanese kimono with seasonal motifs, or European ball gowns with structured silhouettes. Current models often default to the most represented cultural interpretation in their training data, revealing a fundamental gap in their ability to understand implicit cultural contexts and generate culturally-appropriate content without explicit guidance.

To address this cultural bias, efficient fine-tuning approaches have gained attention as a means to adapt pretrained models to specific cultural domains. Low-Rank Adaptation (LoRA) [6, 7] has emerged as a particularly promising solution, enabling customization of large diffusion models with minimal computational overhead. However, when applied to cultural adaptation, conventional LoRA approaches face two critical limitations. First, they apply *uniform adaptations* across all features without considering semantic relevance, leading to inefficient parameter usage and potential interference with unrelated image aspects. Second, they rely solely on reconstruction-based objectives, lacking explicit semantic guidance to ensure cultural authenticity.

While recent attempts to address cultural bias have explored data augmentation and prompt engineering techniques [8], these approaches face significant limitations. They either require extensive manual data collection or depend on carefully crafted prompts, limiting their scalability and practical applicability. More fundamentally, they fail to address the core challenge: enabling models to develop *intrinsic understanding* of cultural semantics that can be activated automatically across diverse generation scenar-

---

[1]Additional qualitative results and visual comparisons are available at our project page: REMOVED-FOR-REVIEW

ios without explicit cultural markers in the input.

Building on these observations, we introduce **K-StyleLoRA**, a novel framework that addresses cultural bias through semantically-guided adaptation. Our approach is motivated by two key insights: (1) pre-trained vision-language models like CLIP [9] possess rich cross-modal knowledge that can guide cultural adaptation, and (2) selective feature enhancement based on cultural relevance can achieve superior authenticity with fewer parameters than uniform adaptation approaches.

Our method integrates two complementary mechanisms into the LoRA adaptation process. *CLIP-Guided Information Gating* leverages CLIP's cross-modal understanding to compute cultural relevance scores for input features, dynamically modulating LoRA adaptations to enhance culturally-relevant aspects while suppressing irrelevant ones. This selective approach enables more efficient adaptation and demonstrates superior few-shot learning capabilities, making it particularly suitable for cultural domains where large-scale training data may be limited. Complementing this, our *Cultural Semantic Loss* provides explicit training guidance by encouraging generated images to align with cultural concepts through CLIP-based similarity optimization, ensuring the model develops genuine cultural understanding rather than superficial pattern matching.

We validate our approach through comprehensive experiments on Korean traditional art generation, a challenging domain requiring sophisticated understanding of aesthetic principles, color palettes, and compositional elements. Our evaluation employs four distinct prompt categories—explicit cultural references, implicit cultural cues, photo-style descriptions, and generic prompts without cultural markers—enabling thorough assessment of the model's cultural understanding and transfer capabilities across different prompt types.

The main contributions of this work are:

- We propose **K-StyleLoRA**, a novel framework that combines CLIP-guided information gating with cultural semantic loss for effective cultural adaptation in diffusion models.

- We introduce *CLIP-Guided Information Gating*, a mechanism that dynamically modulates LoRA adaptations based on cultural relevance scores, enabling targeted parameter updates while preserving general knowledge and demonstrating superior few-shot learning capabilities.

- We design a *Cultural Semantic Loss* that provides explicit semantic guidance through CLIP-based similarity optimization, ensuring authentic cultural representation and enabling effective learning even with limited cultural training data.

Our framework addresses a critical gap in current generative AI systems and provides a principled approach to mitigating cultural bias. While we focus on Korean traditional art as our primary testbed, the methodology is designed for generalizability to diverse cultural domains and generation tasks, opening new avenues for inclusive and culturally-aware generative AI.

## 2. Related Work

### 2.1. Text-to-Image Generation and Cultural Bias

Recent years have witnessed remarkable progress in text-to-image generation, driven primarily by advances in diffusion models [10–14]. Latent Diffusion Models (LDMs) [1] significantly improved computational efficiency by operating in the latent space of pre-trained autoencoders. Building upon this foundation, Stable Diffusion demonstrated remarkable text-to-image generation capabilities [15–17]. More recently, Stable Diffusion XL (SDXL) [18] introduced significant architectural improvements including a larger UNet, refined conditioning mechanisms, and enhanced text encoders, achieving state-of-the-art performance in high-resolution image generation. Other notable models like DALL-E 2 [3] and Imagen [2] have also achieved unprecedented quality in generating photorealistic images from textual descriptions.

Despite these advances, recent studies have revealed significant cultural biases in AI models [4, 5, 19, 20]. These biases manifest as systematic underrepresentation of non-Western cultures and stereotypical portrayals that fail to capture authentic cultural aesthetics. SDXL, while demonstrating superior generation quality, inherits similar cultural biases from its training data, predominantly reflecting Western cultural perspectives. Limited work has explicitly addressed cultural representation in generative models. Early approaches focused on dataset augmentation strategies [21, 22] or bias mitigation through careful training data curation. Others employ prompt engineering techniques [8] or style transfer [23] methods to improve cultural representation. More recent efforts have explored incorporating cultural knowledge through cultural concept embeddings or style codes. However, these methods often require extensive cultural annotations and lack the cross-modal semantic understanding necessary for authentic cultural representation. Our work addresses this critical limitation by developing culturally-aware adaptation techniques that leverage pre-trained vision-language models to enhance SDXL's cultural representation capabilities.

### 2.2. Parameter-Efficient Fine-tuning

The computational cost of full fine-tuning large-scale generative models has motivated research into parameter-efficient adaptation methods [24]. Low-Rank Adaptation (LoRA) [6] emerged as a particularly effective approach, learning low-rank decompositions of weight updates that can be efficiently merged with pre-trained param-

eters. AdaLoRA [25] extends this by adaptively allocating ranks across different layers based on importance scores.

Recent work has explored block-wise adaptation strategies. B-LoRA [26, 27] introduces selective application of LoRA to specific transformer blocks, while Frenkel et al. [28] use B-LoRA for style-content separation in diffusion models. However, these approaches focus on general artistic styles rather than cultural-specific adaptations and lack explicit semantic guidance for cultural authenticity.

In diffusion models, several works have applied LoRA for domain adaptation [29, 30], including applications to SDXL for various customization tasks. However, these methods treat all features uniformly without considering semantic relevance to the target domain. Furthermore, selective adaptation approaches have shown promise for few-shot learning scenarios, where limited training data requires careful parameter allocation to prevent overfitting while maintaining adaptation effectiveness. Our K-StyleLoRA addresses these limitations through CLIP-guided selective feature learning specifically designed for cultural adaptation of SDXL, demonstrating superior performance even with limited cultural training data.

### 2.3. Vision-Language Models for Guidance

Vision-language models, particularly CLIP [9], have demonstrated remarkable capabilities in understanding semantic relationships across modalities. Several works have leveraged CLIP for guiding image generation and editing. CLIP-Guided Diffusion [31] uses CLIP to steer the generation process toward desired text prompts during inference. CLIPStyler [32] employs CLIP for artistic style transfer by optimizing images to match style descriptions.

More recently, attention has turned to using CLIP for training guidance rather than inference-time steering. InstructPix2Pix [33] uses CLIP features to guide instruction-based image editing, while other works explore CLIP-based loss functions for improved text-image alignment during training. However, none of these approaches specifically address cultural adaptation or employ CLIP for selective feature learning in parameter-efficient fine-tuning of large-scale models like SDXL.

## 3. Method

We present K-StyleLoRA, a novel framework that integrates CLIP's cross-modal understanding with Low-Rank Adaptation for culturally-aware image generation. As illustrated in Figure 1, our approach addresses the limitation of conventional LoRA methods that apply uniform adaptations without semantic awareness. Our K-StyleLoRA consists of two key innovations: CLIP-guided information gating for selective feature modulation and cultural semantic loss for global consistency enforcement.
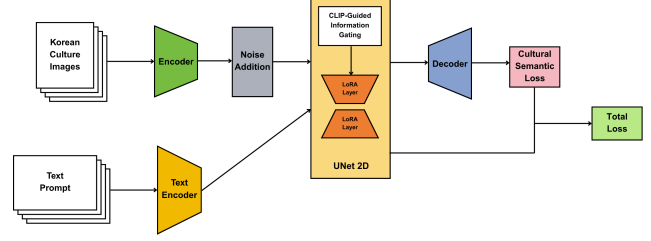


Figure 1. Overview of K-StyleLoRA framework. Our method integrates CLIP-guided information gating with cultural semantic loss for culturally-aware adaptation. (a) CLIP-Guided Information Gating dynamically modulates LoRA adaptations based on cultural relevance scores. (b) Cultural Semantic Loss provides semantic guidance by optimizing CLIP-based similarity between generated images and Korean cultural concepts. (c) The combined approach enables both explicit cultural activation via [v] tokens and implicit cultural transfer for generic prompts.

### 3.1. CLIP-Guided Information Gating

#### 3.1.1. Cultural Relevance Assessment

We employ a shared CLIP ViT-B/32 model to maintain consistent cultural understanding across all components. Korean cultural concepts are pre-encoded and normalized:

$$T_{\text{cultural}} = \{\mathcal{F}_{\text{norm}}(\text{CLIP}_{\text{text}}(c_k))\}_{k=1}^{8} \quad (1)$$

where concepts include "Korean traditional painting", "Korean hanbok clothing", "Korean traditional architecture", and other cultural descriptors.

For input features $h \in \mathbb{R}^{B \times S \times D}$ from UNet attention layers, we compute cultural relevance by projecting features to CLIP's visual space:

$$v = \mathcal{F}_{\text{norm}}(W_{\text{proj}} \cdot \text{mean}(h, \dim = 1)) \quad (2)$$

$$s = \sigma \cdot (v \cdot T_{\text{cultural}}^T) \quad (3)$$

where $W_{\text{proj}} \in \mathbb{R}^{512 \times D}$ projects features to CLIP space and $\sigma$ controls guidance strength.

#### 3.1.2. Selective Feature Modulation

The information gating network generates element-wise modulation weights based on both original features and cultural relevance:

$$g = \text{Sigmoid}(\text{MLP}([h_{\text{flat}}; s_{\text{expanded}}])) \in [0, 1]^{BS \times D} \quad (4)$$

where sigmoid activation ensures gating weights between 0 (suppress) and 1 (enhance). The enhanced LoRA layer then applies selective adaptation:

$$y = h + \frac{\alpha}{r} \cdot \text{reshape}(g \odot \text{flatten}(BAh)) \quad (5)$$

where $B$ and $A$ are LoRA matrices, $\odot$ denotes element-wise multiplication, and the gating mechanism $g$ ensures only culturally-relevant features are adapted.

## 3.2. Cultural Semantic Loss

To enforce global cultural consistency, we introduce a Cultural Semantic Loss that leverages CLIP's vision encoder to assess the cultural alignment of generated images during training.

### 3.2.1. Cultural Guidance Scaling

Following classifier-free guidance principles, we maintain both conditional and unconditional text embeddings:

$$T_{\text{cultural}} = \{\text{CLIP}_{\text{text}}(c_k)\}_{k=1}^{8} \tag{6}$$

$$T_{\text{uncond}} = \text{CLIP}_{\text{text}}(\text{``''}) \tag{7}$$

For generated images $\hat{x}$ decoded from predicted latents during training, we compute CLIP visual features and measure similarities:

$$V = \mathcal{F}_{\text{norm}}(\text{CLIP}_{\text{vision}}(\text{preprocess}(\hat{x}))) \tag{8}$$

$$S_{\text{cultural}} = V \cdot T_{\text{cultural}}^{T} \tag{9}$$

$$w_{\text{concepts}} = \text{Softmax}(S_{\text{cultural}}) \tag{10}$$

$$s_{\text{cond}} = \sum_{k} S_{\text{cultural}} \odot w_{\text{concepts}} \tag{11}$$

$$s_{\text{uncond}} = V \cdot T_{\text{uncond}} \tag{12}$$

We apply cultural guidance scaling similar to CFG:

$$s_{\text{guided}} = s_{\text{uncond}} + \gamma(s_{\text{cond}} - s_{\text{uncond}}) \tag{13}$$

where $\gamma = 7.5$ is the cultural guidance scale. The Cultural Semantic Loss maximizes cultural alignment:

$$\mathcal{L}_{\text{cultural}} = -s_{\text{guided}} \tag{14}$$

## 3.3. Training Objective and Implementation

Our complete training objective combines standard diffusion loss with cultural semantic guidance:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{diffusion}} + \lambda\mathcal{L}_{\text{cultural}} \tag{15}$$

where $\mathcal{L}_{\text{diffusion}} = \mathbb{E}_{t,\epsilon\sim\mathcal{N}(0,1)}\|\epsilon - \epsilon_\theta(z_t, t, c)\|^2$ and $\lambda = 0.1$.

# 4. Experiments

In this section, we present comprehensive experiments to evaluate the effectiveness of K-StyleLoRA for culturally-aware text-to-image generation. We conduct extensive comparisons with baseline methods, ablation studies on key components, and both quantitative and qualitative analyses on Korean traditional art generation.

## 4.1. Experimental Setup

### 4.1.1. Dataset

We curate a dataset of Korean traditional cultural images consisting of 128 high-quality images collected from copyright-free online sources. Training data consists of Korean cultural images from public heritage collections.

### 4.1.2. Implementation Details

We implement K-StyleLoRA on top of Stable Diffusion XL using PyTorch and the Diffusers library. Training is performed on NVIDIA RTX 4080 SUPER. We use a batch size of 1 per device with gradient accumulation steps of 4, resulting in an effective batch size of 4. The resolution is set to 1024×1024 pixels to match SDXL's native resolution.

The shared CLIP model uses ViT-B/32 architecture for computational efficiency. LoRA rank is set to $r = 4$ across all experiments. The cultural loss weight $\lambda$ is set to 0.1, cultural guidance scale $\gamma$ to 7.5, and guidance strength $\sigma$ to 1.0. Training is performed for 10 epochs with AdamW optimizer using learning rate $1 \times 10^{-4}$.

### 4.1.3. Baseline Methods

We compare K-StyleLoRA against several strong baseline methods:

**Vanilla SDXL**: The original Stable Diffusion XL model without any cultural adaptation, serving as the baseline for cultural representation capability.

**Standard LoRA**: Conventional LoRA fine-tuning applied to the same attention blocks as our method but without CLIP guidance or cultural semantic loss ($\sigma = 0$, $\lambda = 0$).

**LoRA + CLIP Loss**: Standard LoRA augmented with cultural semantic loss but without CLIP-guided information gating ($\sigma = 0$, $\lambda = 0.1$).

All baseline methods are trained using the same dataset, computational resources, and training duration to ensure fair comparison.

## 4.2. Evaluation Metrics

### 4.2.1. Quantitative Metrics

We employ automated metrics to comprehensively assess generation quality and cultural authenticity:

**Cultural Similarity Score (CSS)**: We compute the average CLIP similarity between generated images and the 8 pre-defined Korean cultural concepts used in our training. This metric directly measures cultural alignment:

$$\text{CSS} = \frac{1}{N}\sum_{i=1}^{N}\max_{k=1}^{K}\text{sim}(\text{CLIP}(I_i), T_k) \tag{16}$$

where $I_i$ are generated images, $T_k$ are Korean cultural concept embeddings, and $N$ is the number of generated images.

**CLIP Score**: Standard CLIP score between generated images and input prompts to evaluate text-image alignment.

**LPIPS**: We compute average LPIPS distance between pairs of generated images from the same prompt to measure generation diversity.

## 4.3. Main Results

Table 1 presents the quantitative comparison of K-StyleLoRA against baseline methods across key evaluation

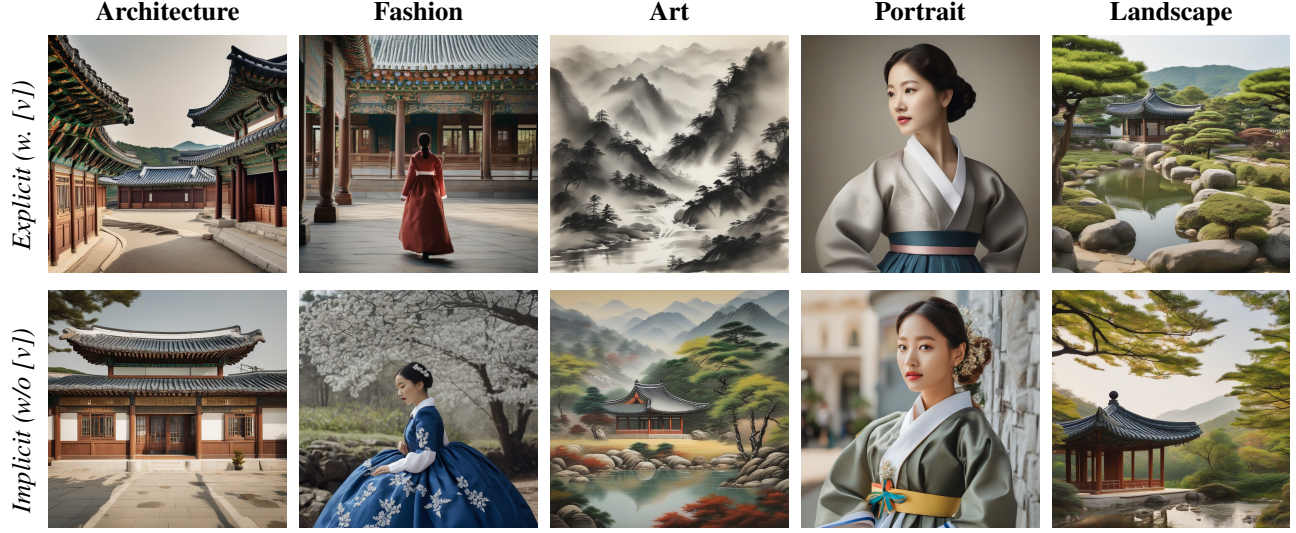| Architecture | Fashion | Art | Portrait | Landscape |
|---|---|---|---|---|



Figure 2. K-StyleLoRA qualitative results organized by Korean cultural categories. Each column represents a different cultural domain. Top row shows explicit cultural activation using [v] token representing "Korean traditional" (e.g., "a woman wearing [v] dress"), while bottom row demonstrates implicit cultural transfer for equivalent prompts without tokens (e.g., "a woman wearing a traditional dress"). Our method successfully applies Korean aesthetics in both scenarios while maintaining high visual quality and prompt adherence.

| Method | CSS(All) ↑ | CSS(Exp) ↑ | CSS(Gen) ↑ | CLIP ↑ |
|---|---|---|---|---|
| Vanilla SDXL | 0.291±0.04 | 0.299±0.02 | 0.250±0.04 | 0.324±0.02 |
| Standard LoRA | 0.290±0.04 | 0.295±0.03 | 0.255±0.05 | 0.323±0.02 |
| LoRA + CLIP | 0.293±0.04 | 0.302±0.03 | 0.261±0.04 | 0.327±0.02 |
| **K-StyleLoRA** | **0.298±0.04** | **0.309±0.03** | **0.274±0.03** | **0.327±0.02** |

Table 1. Quantitative evaluation results on Korean traditional art generation.

metrics. We evaluate cultural alignment using our Cultural Similarity Score (CSS) across different prompt categories: explicit cultural prompts (CSS-Exp), and generic prompts requiring implicit cultural transfer (CSS-Gen).

K-StyleLoRA demonstrates improved performance across cultural similarity metrics, achieving the highest CSS(Exp) score of 0.309, representing a 9.6% improvement over vanilla SDXL on generic prompts (CSS-Gen): 0.274 vs 0.250. Our method also maintains superior text-image alignment while preserving cultural authenticity.

### 4.4. Ablation Studies

| Method | CSS(All) ↑ | CSS(Gen) ↑ | CLIP ↑ | LPIPS ↑ |
|---|---|---|---|---|
| Standard LoRA | 0.290±0.04 | 0.255±0.05 | 0.323±0.02 | 0.725±0.06 |
| + CLIP Loss Only | 0.293±0.04 | 0.261±0.04 | 0.327±0.02 | 0.755±0.07 |
| + CLIP Gating Only | 0.293±0.04 | 0.262±0.04 | 0.325±0.02 | 0.724±0.06 |
| K-StyleLoRA ($\sigma = 0.5$) | 0.289±0.04 | 0.250±0.03 | 0.324±0.02 | 0.761±0.08 |
| K-StyleLoRA ($\sigma = 1.0$) | 0.292±0.04 | 0.257±0.04 | **0.330±0.02** | 0.741±0.06 |
| K-StyleLoRA ($\sigma = 1.5$) | **0.298±0.04** | **0.274±0.03** | 0.327±0.02 | 0.749±0.06 |

Table 2. Ablation Study on K-StyleLoRA Components and Guidance Strength.

To understand the contribution of each component, we conduct comprehensive ablation studies examining our two key innovations: CLIP-guided information gating and cultural semantic loss.

The ablation study reveals that both CLIP loss and CLIP gating contribute similarly to cultural representation. The guidance strength analysis shows that stronger CLIP guidance yields the best cultural similarity scores, demonstrating the effectiveness of our CLIP-guided information gating approach.

### 4.5. Qualitative Results

Figure 2 presents a comprehensive showcase of K-StyleLoRA's generation capabilities across different prompt categories, demonstrating the method's ability to maintain cultural authenticity while producing diverse and high-quality outputs.

Figure 3 demonstrates the comparison between vanilla SDXL and K-StyleLoRA, showing how our method applies Korean cultural aesthetics to various prompt types.

**Cultural Authenticity**: K-StyleLoRA consistently generates images with authentic Korean cultural elements across all prompt categories, accurately capturing traditional Korean aesthetics including proper proportions, color palettes, and cultural motifs.

**Implicit Cultural Transfer**: The method demonstrates strong cultural transfer capabilities even when cultural elements are not explicitly mentioned in the prompt, showing Korean-influenced design elements while maintaining prompt adherence.

**Architecture**

**Portrait**



Figure 3. Direct comparison between Vanilla SDXL (top) and K-StyleLoRA (bottom) demonstrating implicit cultural transfer capabilities. Architecture examples use generic prompts like "traditional building" and "traditional architecture". Portrait examples use prompts such as "a woman wearing a traditional dress" and "a woman in silk dress". K-StyleLoRA automatically applies Korean cultural elements without explicit Korean cultural keywords, producing culturally-specific outputs while maintaining high visual quality.

## 4.6. Limitations

While K-StyleLoRA significantly improves cultural representation, several limitations remain. First, our method is inherently dependent on CLIP's performance and cultural understanding capabilities. If CLIP has biases or limitations in recognizing certain cultural elements, these may affect our generation system's performance. Future work could explore using more diverse vision-language models or developing cultural-specific evaluation metrics to better assess cross-cultural generation quality. Second, the method requires careful curation of cultural concept embeddings, which may not capture all nuances of a cultural domain. Very abstract or conceptual cultural elements may still be challenging to represent accurately. Adaptive cultural embedding techniques could address this limitation. Third, our experiments are conducted on a limited dataset and specific experimental settings. More extensive evaluation across diverse cultural domains, larger datasets, and varied experimental configurations would be necessary to fully validate the generalizability and robustness of our approach. Scaling to larger datasets and cross-cultural validation studies would strengthen these findings.

## 5. Conclusion

In this paper, we presented K-StyleLoRA, a novel framework for culturally-aware image generation that addresses cultural bias in existing diffusion models. Our approach introduces two key innovations: CLIP-Guided Information Gating for selective feature learning based on cultural relevance, and Cultural Semantic Loss for global semantic guidance through CLIP-based similarity optimization. Extensive experiments on Korean traditional art demonstrate that K-StyleLoRA achieves a 9.6% improvement in cultural similarity over vanilla SDXL while maintaining text-image alignment and generation diversity. Our framework establishes semantic-aware adaptation as a powerful paradigm for cultural representation, offering a scalable approach that can be extended to diverse cultural contexts and generation tasks. This work opens new directions for parameter-efficient cultural adaptation, promoting algorithmic fairness and cultural diversity in generative AI.

## References

[1] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022. 1, 2

[2] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans, *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *Advances in Neural Information Processing Systems*, vol. 35, pp. 36479–36494, 2022. 2

[3] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, "Hierarchical text-conditional image generation with clip latents," *arXiv preprint arXiv:2204.06125*, 2022. 1, 2

[4] F. Bianchi, P. Kalluri, E. Durmus, F. Ladhak, M. Cheng, D. Nozza, T. Hashimoto, D. Jurafsky, J. Zou, and A. Caliskan, "Easily accessible text-to-image generation amplifies demographic stereotypes at large scale," *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pp. 1493–1504, 2023. 1, 2

[5] J. Cho, A. Zala, and M. Bansal, "Dall-e 2 is seeing double: Flaws in word-to-image generation models," *Proceedings of*

the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pp. 1–14, 2023. 1, 2

[6] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," *International Conference on Learning Representations*, 2022. 1, 2

[7] N. Ruiz, Y. Li, V. Jampani, W. Wei, T. Hou, Y. Pritch, N. Wadhwa, M. Rubinstein, and K. Aberman, "Hyperdreambooth: Hypernetworks for fast personalization of text-to-image models," 2023. 1

[8] J. Liu, C. Wang, S. Joty, C. Xiong, and R. Socher, "Cultural prompting: Activating cultural knowledge in language models for bias reduction," *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, pp. 4556–4571, 2023. 1, 2

[9] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, *et al.*, "Learning transferable visual models from natural language supervision," *International Conference on Machine Learning*, pp. 8748–8763, 2021. 2, 3

[10] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020. 2

[11] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," *International Conference on Learning Representations*, 2021.

[12] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen, "Glide: Towards photorealistic image generation and editing with text-guided diffusion models," in *Proceedings of the 39th International Conference on Machine Learning (ICML)*, pp. 16784–16804, 2022.

[13] J. Li, R. Selvaraju, A. Gotmare, S. Joty, C. Xiong, and S. Hoi, "Align before fuse: Vision and language representation learning with momentum distillation," in *Advances in Neural Information Processing Systems*, vol. 34, pp. 9694–9705, 2021.

[14] O. Patashnik, Z. Wu, E. Shechtman, D. Cohen-Or, and D. Lischinski, "Styleclip: Text-driven manipulation of stylegan imagery," *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2085–2094, 2021. 2

[15] L. Zhang, A. Rao, and M. Agrawala, "Adding conditional control to text-to-image diffusion models," *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3836–3847, 2023. 2

[16] M. Li, T. Yang, H. Kuang, J. Wu, Z. Wang, X. Xiao, and C. Chen, "Controlnet++: Improving conditional controls with efficient training," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2024.

[17] C. Mou, X. Wang, L. Xie, Y. Wu, J. Zhang, Z. Qi, Y. Shan, and X. Qie, "T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, vol. 38, pp. 4296–4304, 2024. 2

[18] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach, "Sdxl: Improving latent diffusion models for high-resolution image synthesis," in *Proceedings of the 12th International Conference on Learning Representations (ICLR)*, 2024. 2

[19] S. L. Blodgett, S. Barocas, H. Daumé III, and H. Wallach, "Language (technology) is power: A critical survey of "bias" in nlp," *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5454–5476, 2020. 2

[20] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, "On the dangers of stochastic parrots: Can language models be too big?," *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pp. 610–623, 2021. 2

[21] J. An *et al.*, "Adaptive augmentation for effectively mitigating dataset bias," in *Asian Conference on Computer Vision (ACCV)*, 2022. 2

[22] S. Sharma *et al.*, "Data augmentation for discrimination prevention and bias disambiguation," in *ACM Conference on Fairness, Accountability, and Transparency (FAccT)*, 2020. 2

[23] L. A. Gatys, A. S. Ecker, and M. Bethge, "Image style transfer using convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2414–2423, 2016. 2

[24] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, "Qlora: Efficient finetuning of quantized llms," in *Advances in Neural Information Processing Systems*, vol. 36, 2023. 2

[25] Q. Zhang, M. Chen, A. Bukharin, P. He, Y. Cheng, W. Chen, and T. Zhao, "Adaptive budget allocation for parameter-efficient fine-tuning," *International Conference on Learning Representations*, 2023. 3

[26] Y. Wang, Z. Chen, Q. Liu, M. Zhang, and J. Li, "B-lora: Block-wise low-rank adaptation for efficient fine-tuning of large language models," *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 5847–5859, 2023. 3

[27] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, "Block-wise lora: Switchable and efficient fine-tuning for large foundation models," *International Conference on Machine Learning*, pp. 21540–21558, 2023. 3

[28] Y. Frenkel, Y. Vinker, A. Shamir, and D. Cohen-Or, "Implicit style-content separation using b-lora," in *ECCV*, pp. 181–198, Springer, 2024. 3

[29] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman, "Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22500–22510, 2023. 3

[30] R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-Or, "An image is worth one word: Personalizing text-to-image generation using textual inversion," *European Conference on Computer Vision*, pp. 208–224, 2022. 3

[31] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," *Advances in Neural Information Processing Systems*, vol. 34, pp. 8780–8794, 2021. 3

[32] G. Kwon and J. C. Ye, "Clipstyler: Image style transfer with a single text condition," *Proceedings of the IEEE/CVF*

7

*Conference on Computer Vision and Pattern Recognition*, pp. 17292–17301, 2022. 3

[33] T. Brooks, A. Holynski, and A. A. Efros, "Instructpix2pix: Learning to follow image editing instructions," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18392–18402, 2023. 3